

Europäisches Patentamt  
European Patent Office  
Office européen des brevets



(11) **EP 0 459 232 B1**

A1

(12) **EUROPEAN PATENT SPECIFICATION**

(45) Date of publication and mention  
of the grant of the patent:  
**09.12.1998 Bulletin 1998/50**

(51) Int Cl.<sup>6</sup>: **G06F 9/38**

(21) Application number: **91107898.8**

(22) Date of filing: **16.05.1991**

(54) **Partially decoded instruction cache and method therefor**

Cache-Speicher von partiell decodierten Befehlen und Verfahren hierfür

Antémémoire d'instructions décodées partiellement et méthode correspondante

(84) Designated Contracting States:  
**DE FR GB IT**

(30) Priority: **29.05.1990 US 529869**

(43) Date of publication of application:  
**04.12.1991 Bulletin 1991/49**

(73) Proprietor: **NATIONAL SEMICONDUCTOR  
CORPORATION**  
**Santa Clara California 95051-8090 (US)**

(72) Inventors:  
• **Alpert, Donald B.**  
**Herzli (IL)**  
• **Avnon, Dror**  
**Netanya (IL)**

• **Ben-Meir, Amos**  
**Ramat Aviv (IL)**  
• **Talmudi, Ran**  
**Raanana (IL)**

(74) Representative: **Sparing Röhl Henseler**  
**Patentanwälte**  
**Postfach 14 04 43**  
**40074 Düsseldorf (DE)**

(56) References cited:  
**EP-A- 0 363 222** **WO-A-90/03001**  
**US-A- 4 873 629**

• **MICROPROCESSORS AND MICROSYSTEMS,**  
**vol. 13, no. 9, November 1989, LONDON, GB;**  
**pages 579 - 587; G. B. STEVEN ET AL: 'HARP: A**  
**Parallel Pipelined RISC Processor'**

Note: Within nine months from the publication of the mention of the grant of the European patent, any person may give notice to the European Patent Office of opposition to the European patent granted. Notice of opposition shall be filed in a written reasoned statement. It shall not be deemed to have been filed until the opposition fee has been paid. (Art. 99(1) European Patent Convention).

## Description

The present invention relates to microprocessor architectures and, in particular, to a microprocessor that partially decodes instructions retrieved from external memory before storing them in an internal instruction cache. Partially decoded instructions are retrieved from the internal cache for either parallel or sequential execution by multiple, parallel, pipelined functional units.

In recent years, there has been a trend in the design of microprocessor architectures from Complex Instruction Set Computers (CISC) toward Reduced Instruction Set Computers (RISC) to achieve high performance while maintaining simplicity of design.

In a CISC architecture, each macroinstruction received by the processor must be decoded internally into a series of microinstruction subroutines. These microinstruction subroutines are then executed by the microprocessor.

In a RISC architecture, the number of macroinstructions which the processor can understand and execute is greatly reduced. Further, those macroinstructions which the processor can understand and execute are very basic so that the processor either does not decode them into any microinstructions (the macroinstruction is executed in its macro form) or the decoded microinstruction subroutine involves very few microinstructions.

The transition from CISC architectures to RISC architectures has been driven by two fundamental developments in computer design that are now being extensively applied to microprocessors. These developments are integrated cache memory and optimizing compilers.

A cache memory is a small, high speed buffer located between the processor and main memory to hold the instructions and data most recently used by the processor. Experience shows that computers very commonly exhibit strong characteristics of locality in their memory references. That is, references tend to occur frequently either to locations that have recently been referred to (temporal locality) or to locations that are near others that have recently been referred to (spatial locality). As a consequence of this locality, a cache memory that is much smaller than main memory can capture the large majority of a program's memory references. Because the cache memory is relatively small, it can be realized from a faster memory technology than would be economical for the much larger main memory.

Before the development of cache memory techniques for use in mainframe computers, there was a large imbalance between the cycle time of a processor and that of memory. This imbalance was a result of the processor being realized from relatively high speed bipolar semiconductor technology and the memory being realized from much slower magnetic-core technology. The inherent speed difference between logic and memory spurred the development of complex instruction sets that would permit the fetching of a single instruction from memory to control the operation of the processor for

several clock cycles. The imbalance between processor and memory speeds was also characteristic of the early generations of 32-bit microprocessors. Those microprocessors would commonly take 4 or 5 clock cycles for each memory access.

Without the introduction of integrated cache memory, it is unlikely that RISC architectures would have become competitive with CISC architectures. Because a RISC processor executes more instructions than does a CISC processor to accomplish the same task, a RISC processor can deliver performance equivalent to that of a CISC only if a faster and more expensive memory system is employed. Integrated cache memory enables a RISC processor to fetch an instruction in the same time required to execute the instruction by an efficient processor pipeline.

The second development that has led to the effectiveness of RISC architectures is optimizing compilers. A compiler, which may be implemented in either hardware or software, translates a computer program from the high-level language used by the programmer into the machine language understood by the computer.

For many years after the introduction of high-level languages, computers were still extensively programmed in assembly language. Assembly language is a low-level source code language employing crude mnemonics that are more easily remembered by the programmer than object-code or binary equivalents. The advantages of improved software productivity and translatability of high-level language programming were clear, but simple compilers produced inefficient code. Early generations of 32-bit microprocessors were developed with consideration for assembly language programming and simple compilers.

More recently, advances in compiler technology are being applied to microprocessors. Optimizing compilers can analyze a program to allocate large numbers of registers efficiently and to manage processor pipeline resources. As a consequence, high-level language programs can execute with performance comparable to or exceeding that of assembly programs.

Many of the leading pioneers in RISC developments have been compiler specialists who have demonstrated that optimizing compilers can produce highly efficient code for simple, regular architectures.

Highly integrated single-chip microprocessors employ both pipelined and parallel execution to improve performance. Pipelined execution means that while the microprocessor is fetching one instruction, it can be simultaneously decoding a second instruction, reading source operands for a third instruction, calculating results for a fourth instruction and writing results for a fifth instruction. Parallel execution means that the microprocessor can initiate the operands for two or more independent instructions simultaneously in separate functional units.

As stated above, one of the main challenges in designing a high-performance microprocessor with multi-

ple, pipelined functional units is to provide sufficient instruction memory on-chip and to access the instruction memory efficiently to control the functional units.

The requirement for efficient control of a microprocessor's functional units dictates a regular instruction format that is simple to decode. However, in conventional microprocessor architectures, instructions in main memory are highly encoded and of variable length to make efficient use of space in main memory and the limited bandwidth available between the microprocessor and the main memory.

The present invention is defined by the independent claims 1 and 8 and provides a processor and corresponding method that resolves the conflicting requirements for efficient use of main memory storage space and efficient control of the functional units by partially decoding instructions retrieved from main memory before placing them into the microprocessor's integrated instruction cache. Preferably, each entry in the instruction cache has two slots for partially decoded instructions. One slot controls one of the microprocessor's execution pipelines and a port to its data cache. The second slot controls a second execution pipeline, or one of the microprocessor's floating point units, or a control transfer instruction. An instruction decoding unit, or loader, decodes instructions from their compact format as stored in main memory and places them into the two slots of the instruction cache entry according to their functions. Auxiliary information may also be placed in the cache entry along with the instruction to control parallel execution and emulation of complex instructions. A bit in each cache entry may indicate whether the instructions in the two slots for that entry are independent, so that they can be executed in parallel, or dependent, so that they must be executed sequentially. Using a single bit for this purpose allows two dependent instructions to be stored in the slots of a single cache entry. Otherwise, the two instructions would have to be stored in separate entries and only one-half of the cache memory would be utilized in those two entries.

Some features of the independent claims are known per se.

US-A-4,873,629 discloses a computer configured for optimizing the processing rate of instructions and a corresponding method. The computer includes a main memory, a cache unit, and a central processing unit. According to this document (encoded) instructions retrieved from the main memory are "cracked", i.e. the address fields of the instructions are decoded such that they can be stored in a logical instruction cache unit. When the cracked instructions are retrieved from the cache unit for subsequent execution, they are (sequentially) sent to an output buffer and decoder means where a decoding step occurs and furthermore decoded program count and displacement information is generated. Decoded instructions then are sent to ALUs for execution thereof.

The technique according to this document therefore

requires time-consuming multiple storage of the instructions, and especially teaches to decouple the cache means and the functional units by intermediate buffer means.

Another technique is disclosed by EP-A-O 363 222. This document relates to apparatus and method for concurrent dispatch of instruction words. These instructions are capable to be separately and substantially simultaneously received by distinct functional units, i.e. a floating point unit and an integer unit being part of a processor. Instructions coming from an external source are retrieved and stored in an instruction cache, wherein the cache is partitioned into even and odd cache sections, which each are subsequently connected via logic means to functional units. The logic units include decoding means that are capable to decode the encoded instructions after reception. Accordingly, the instructions cached in the instruction cache are still encoded.

The known architecture is suitable for Complex Instruction Set Computers (CISC) with a central processing unit, but will encounter problems when being transferred to a Reduced Instruction Set Computer (RISC) architecture. A problem arising when processing instructions in a RISC architecture is that the number of instructions understood by the processor is greatly reduced, at the same time, these instructions are performed very quickly. The operation of completely decoding a retrieved encoded instruction bears an unpredictable doubt concerning required time when performing the decoding operation so that processor time is difficult to optimize.

An article by Stevens, G.B. et al, "HARP": A parallel pipelined RISC processor", Microprocessor and Microsystems Vol. 13, No. 9, November 1989, pp 579-587, London, GB, relates to a compiler that packs independent "HARP" instructions being executable in parallel into long instruction words. The long instruction words are retrieved from an instruction cache for passing the component short instructions of the long instruction word through a parallel pipeline structure.

There is no teaching that could give a hint to partially decode encoded instructions before storing them in a cache.

WO-A-90 03 001 discloses a CISC system in which encoded instructions are retrieved from the main memory, subsequently stored in a cache means. The instructions, when retrieved from the cache, are partially decoded and stored into an FIFO instruction buffer before execution.

Such FIFO queue inter alia provides gradual decoupling of cache and processor unit for buffering time runouts due to partially decoding. It is evident that this solution provides more complicated architecture, and of course does not teach the invention.

A better understanding of the features and advantages of the present invention will be obtained by reference to the following detailed description of the invention and accompanying drawings which set forth an il-

illustrative embodiment in which the principles of the invention are utilized.

Figure 1 is a block diagram illustrating a microprocessor architecture that incorporates the concepts of the present invention.

Figure 2 is a block diagram illustrating the structure of a partially decoded instruction cache utilized in the figure 1 architecture.

Figure 3 is a simplified representation of a partially decoded entry stored in the instruction cache shown in figure 2.

Figure 4 is a block diagram illustrating the structure of the integer pipelines utilized in the microprocessor architecture shown in figure 1.

Figure 1 shows a block diagram of a microprocessor 10 that includes multiple, pipelined functional units that are capable of executing two instructions in parallel.

The microprocessor 10 includes three main sections: an instruction processor 12, an execution processor 14 and a bus interface processor 16.

The instruction processor 12 includes three modules: an instruction loader 18, an instruction emulator 20 and an instruction cache 22. These modules load instructions from the external system through the bus interface processor 16, store the instructions in the instruction cache 22 and provide pairs of instructions to the execution processor 14 for execution.

The execution processor 14 includes two 4-stage pipelined integer execution units 24 and 26, a double-precision 5-stage pipelined floating point execution unit 28, and a 1024 byte data cache 30. A set of integer registers 32 services the two integer units 24 and 26; similarly, a set of floating point registers 34 services the floating point execution unit 28.

The bus interface processor 16 includes a bus interface unit 36 and a number of system modules 38. The bus interface unit 36 controls the bus accesses requested by both the instruction processor 12 and the execution processor 14. In the illustrated embodiment, the system modules 38 include a timer 40, a direct memory access (DMA) controller 42, an interrupt control unit (ICU) 44 and I/O buffers 46.

As described in greater detail below, the instruction loader 18 partially decodes instructions retrieved from main memory and places the partially decoded instructions in the instruction cache 22. That is, the instruction loader 18 translates an instruction stored in main memory (not shown) into the decoded format of the instruction cache 22. As will also be described in greater detail below, the instruction loader 18 is also responsible for checking whether any dependencies exist between consecutive instructions that are paired in a single instruction cache entry.

The instruction cache 22 contains 512 entries for partially-decoded instructions.

In accordance with one aspect of the present invention, and as explained in greater detail below, each entry in the instruction cache 22 contains either one or two

instructions stored in a partially-decoded format for efficient control of the various functional units of the microprocessor 10.

In accordance with another aspect of the present invention, each entry in instruction cache 22 also contains auxiliary information that indicates whether the two instructions stored in that entry are independent, so that they can be executed in parallel, or dependent, so that they must be executed sequentially.

The instruction emulator 20 executes special instructions defined in the instruction set of the microprocessor 10. When the instruction loader 18 encounters such an instruction, it transfers control to the emulator 20. The emulator is responsible for generating a sequence of core instructions (defined below) that perform the function of a single complex instruction (defined below). In this regard, the emulator 20 provides ROM-resident microcode. The emulator 20 also controls exception processing and self-test operations.

The two 4-stage integer pipelines 24 and 26 perform basic arithmetic/logical operations and data memory references. Each integer pipeline 24,26 can execute instructions at a throughput of one per system clock cycle.

The floating point execution unit 28 includes three sub-units that perform single-precision and double-precision operations. An FPU adder sub-unit 28a is responsible for add and convert operations, a second sub-unit 28b is responsible for multiply operations and a third sub-unit 28c is responsible for divide operations.

When add and multiply operations are alternately executed, the floating point execution unit 28 can execute instructions at a throughput of one instruction per system clock cycle.

Memory references for the floating point execution unit 28 are controlled by one of the integer pipelines 24,26 and can be performed in parallel to floating-point operations.

Data memory references are performed using the 1-Kbyte data cache 30. The data cache 30 provides fast on-chip access to frequently used data. In the event that data are not located in the data cache 30, then off-chip references are performed by the bus interface unit (BIU) 36 using the pipelined system bus 48.

The data cache 30 employs a load scheduling technique so that it does not necessarily stall on misses. This means that the two execution pipelines 24,26 can continue processing instructions and initiating additional memory references while data is being read from main memory.

The bus interface unit 36 can receive requests for main memory accesses from either the instruction processor 12 or the execution processor 14. These requests are sent to the external pipelined bus 48. The external bus can be programmed to operate at half the frequency of the microprocessor 10; this allows for a simple instruction interface at a relatively low frequency while the microprocessor 10 executes a pair of instructions at full

rate.

The instruction set of the microprocessor 10 is partitioned into a core part and a non-core part. The core part of the instruction set consists of performance critical instructions and addressing modes, together with some special-function instructions for essential system operations. The non-core part consists of the remainder of the instruction set. Performance critical instructions and addressing modes were selected based on an analysis and evaluation of the operating system (UNIX in this case) workload and various engineering, scientific and embedded controller applications. These instructions are executed directly as part of the RISC architecture of microprocessor 10.

As stated above, special-function and non-core instructions are emulated in microprocessor 10 by macroinstruction subroutines using sequences of core instructions. That is, instructions that are a part of the overall instruction set of the microprocessor 10 architecture, but that lie outside the directly-implemented RISC core, are executed under control of the instruction emulator 20. When the instruction loader 18 encounters a non-core instruction, it either translates it into a pair of core instructions (for simple instructions like MOV B 1 (R0),0(R1)) or transfers control to the instruction emulator 20. The instruction emulator 20 is responsible for generating a sequence of core instructions that perform the function of the single, complex instruction.

Fig. 2 shows the structure of the instruction cache 22. The instruction cache 22 utilizes a 2-way, set-associative organization with 512 entries for partially decoded instructions. This means that for each memory address there are two entries in the instruction cache 22 where the instruction located at that address can be placed. The two entries are called a "set".

As shown in Fig. 3, each instruction cache entry includes two slots, i.e. Slot A and Slot B. Thus, each entry can contain one or two partially-decoded instructions that are represented with fixed fields for opcode (Opc), source and destination register numbers (R1 and R2, respectively), and immediate values (32b IMM). The entry also includes auxiliary information used to control the sequence of instruction execution, including a bit P that indicates whether the entry contains two consecutive instructions that can be executed in parallel and a bit G that indicates whether the entry is for a complex instruction that is emulated, and additional information representing the length of the instruction(s) in a form that allows fast calculation of the next instruction's address.

Referring back to Fig. 2, associated with each entry in the instruction cache 22 is a 26-bit tag, TAG0 and TAG1, respectively, that holds the 22 most-significant bits, 3 least-significant bits and a User/Supervisor bit of the virtual address of the instruction stored in the entry. In the event that two consecutive instructions are paired in an entry, the tag corresponds to the instruction at the lower address. Associated with the tag are 2 bits that indicate whether the entry is valid and whether it is

locked. For each set there is an additional single bit that indicates the entry within the set that is next to be replaced in a Least-Recently-Used (LRU) order.

The instruction cache 22 is enabled for an instruction fetch if a corresponding bit of the configuration register of microprocessor 10 which is used to enable or disable various operating modes of the microprocessor 10, is 1 and either address translation is disabled or the CI-bit is 0 in the level-2 Page Table Entry (PTE) used to translate the virtual address of the instruction.

If the instruction cache 22 is disabled, then the instruction fetch bypasses the instruction cache 22 and the contents of the instruction cache 22 are unaffected. The instruction is read directly from main memory, partially decoded by the instruction loader 18 to form an entry (which may contain two partially decoded instructions), and transferred to the integer pipelines 24, 26 via the IL BYPASS line for execution.

As shown in Fig. 2, if the instruction cache 22 is enabled for an instruction fetch, then eight bits, i.e. bits PC (10:3), of the instruction's address provided by the program counter (PC) are decoded to select the set of entries where the instruction may be stored. The selected set of four entries is read and the associated tags are compared with the 23 most-significant bits, i.e. PC(31:10), and 2 least-significant bits PC(1:0) of the instruction's virtual address. If one of the tags matches and the matching entry is valid, then the entry is selected for transfer to the integer pipelines 24, 26 for execution. Otherwise, the missing instruction is read directly from main memory and partially decoded, as explained below.

If the referenced instruction is missing from the instruction cache 22 and the contents of the selected set are all locked, then the handling of the reference is identical to that described above for the case when the instruction cache 22 is disabled.

If the referenced instruction is missing from the instruction cache 22 and at least one of the entries in the selected set is not locked, then the following actions are taken. One of the entries is selected for replacement according to the least recently used (LRU) replacement algorithm and then the LRU pointer is updated. If the entry selected for replacement is locked, then the handling of the reference is identical to that described above for the case when the instruction cache 22 is disabled. Otherwise, the missing instruction is read directly from external memory and then partially decoded by instruction loader 18 to form an entry (that may contain two partially decoded instructions) which is transferred to the integer pipelines 24, 26 for execution. If CIIN is not active during the bus cycles to read the missing instruction, then the partially decoded instruction is also written into the instruction cache entry selected for replacement, the associated valid bit is set, and the entry is locked if Lock-Instruction-Cache bit CFG.LIC in the configuration register is 1.

After the microprocessor 10 has completed fetching a missing instruction from external main memory, it will

continue prefetching sequential instructions. For subsequent sequential instruction fetches, the microprocessor 10 searches the instruction cache 22 to determine whether the instruction is located on-chip. If the search is successful or a non-sequential instruction fetch occurs, then the microprocessor 10 ceases prefetching. Otherwise, the prefetched instructions are rapidly available for decoding and executing. The microprocessor 10 initiates prefetches only during bus cycles that would otherwise be idle because no off-chip data references are required.

It is possible to fetch an instruction and lock it into the instruction cache 22 without having to execute the instruction. This can be accomplished by enabling a Debug Trap (DBG) for a Program Counter value that matches two instruction's address. Debug Trap is a service routine that performs actions appropriate to this type of exception. At the conclusion of the DBG routine, the RETurn to Execution (RETX) instruction is executed to resume executing instructions at the point where the exception was recognized. The instruction will be fetched and placed into the Instruction Cache 32 before the trap is processed.

When the instruction which is locked in the instruction cache 22 gets to execution and a Debug Trap on that instruction is enabled, instead of executing the instruction, the processor will jump to the Debug Trap service routine. The service routine may set a breakpoint for the next instruction so that when the processor returns from the service routine, it will not execute the next instruction but rather will go again to the Debug Trap routine.

The process described above, which usually gets executed during system bootstrap, allows the user to store routines in the instruction cache 22, lock them and have them ready for operation without executing them during the locking process.

Further information relating to the architecture of microprocessor 10 and its cache locking capabilities is provided in EP-A-0 459 233.

The contents of the instruction cache 22 can be invalidated by software or by hardware.

The instruction cache 22 is invalidated by software as follows: The entire instruction cache contents, including locked entries, are invalidated while bit CFG.IC of the Configuration Register is 0. The LRU replacement information is also initialized to 0 while bit CFG.IC is 0. Cache Invalidate CINV instruction can be executed to invalidate the entire instruction cache contents. Executing CINV invalidates either the entire cache or only unlocked lines according the instruction's L-option.

The entire instruction cache 22 is invalidated in hardware by activating an  $\overline{\text{INVIC}}$  input signal.

Fig. 3 shows a simplified view of a partially decoded entry stored in the instruction cache 22. As shown in Fig. 3, each entry has two slots for instructions. Slot A controls integer pipeline 24 and the port to data cache 30. Slot B controls the second integer pipe 26, or one of the

floating point units or a control transfer instruction. Slot B can also control the port to data cache 30, but only if slot A is not using the data cache 30. As stated above, instruction loader 18 retrieves encoded instructions from their compact format in main memory and places them into slots A and B according to their functions.

Thus, in accordance with the present invention, the novel aspects of instruction cache 22 include (1) partially decoding instructions for storage in cache memory, (2) placing of instructions into two cache slots according to their function and (3) placing auxiliary information in the cache entries along with the instructions to control parallel execution and emulation of complex instructions.

As further shown in Fig. 3, a bit P in each instruction cache entry indicates whether the instructions in slots A and B are independent, so they can be executed in parallel, or dependent, so they must be executed sequentially.

An example of independent instructions that can be executed in parallel is:

Load 4(R0),R1 ; Addcd 4,R0

An example of dependent instructions requiring sequential execution is:

Add R0, R1 ; Addd R1,R2

Using a single bit for this purpose allows two dependent instructions to be stored in the slots of a single cache entry, otherwise, the two instructions would have to be stored in separate entries and only 1/2 of the instruction cache 22 would be utilized in those two entries.

Fig. 3 also shows a bit G in each instruction cache entry that indicates whether the instructions in slots A and B are emulating a single, more complex instruction from main memory. For example, the loader translates the single instruction ADDD O(R0), R1 into the following pair of instructions in slots A and B and sets the sequential and emulation flags in the entry:

Load 0(R0), Temp

ADDD Temp, R1

In accordance with the pipelined organization of the microprocessor 10, every instruction executed by the microprocessor 10 goes through a series of stages. The two integer pipelines 24, 26 (Fig. 1) are able to work in parallel on instructions pairs. Integer unit 24 and integer unit 26 are not identical, the instructions that can be executed in integer unit 24 being a sub-set of those that can be executed in integer unit 26.

As stated above, instruction fetching is performed by the instruction loader 18 which stores decoded instructions in the instruction cache 22. The integer dual-pipe receives decoded instruction-pairs for execution.

Referring again to Fig. 3, as stated above, an instruction pair consists of two slots: Slot A and Slot B. The instruction in Slot A is scheduled for integer unit 24; the instruction in Slot B is scheduled for integer unit 26. Two instructions belonging to the same pair advance at the same time from one stage of the integer pipeline to the next, except in the case when the instruction in Slot

B is delayed in the instruction decode stage of the pipeline as described below. In this case, the instruction in integer pipeline 24 can advance to the following pipeline stages. However, new instructions cannot enter the pipeline until the instruction decode stage is free in both pipeline unit 24 and pipeline unit 26.

Although the unit 24 and unit 26 instructions are executed in parallel (except in the case of the stall ID-B instruction), the Slot A instruction always logically precedes the corresponding Slot B instruction and, if the Slot A instruction cannot be completed due to an exception, then the corresponding Slot B instruction is discarded.

Referring to Fig. 4, each of the integer pipeline units 24, 26 includes four stages: an instruction decode stage (ID), an execute stage (EX), a memory access stage (ME) and a store result stage (ST).

An instruction is fed into the ID stage of the integer unit for which it is scheduled where its decoding is completed and register source operands are read. In the EX stage, the arithmetic/logical unit of the microprocessor 10 is activated to compute the instruction's results or to compute the effective memory address for Load/Store instructions. In the ME stage, the data cache 30 (Fig. 1) is accessed by Load/Store instructions and exception conditions are checked. In the ST stage, results are written to the register file, or to the data cache 30 in the case of a Store instruction, and Program Status Register (PSR) flags are updated. At this stage, the instruction can no longer be undone.

As further shown in Fig. 4, results from the EX stage and the ME stage can be fed back to the ID stage, thus enabling instruction latency of 1 or 2 cycles.

In the absence of any delays, the dual execution pipeline of microprocessor 10 accepts a new instruction pair every clock cycle (i.e., peak throughput of two instructions per cycle) and scrolls all other instructions down one stage along the pipeline. The dual pipeline includes a global stalling mechanism by which any functional unit can stall the pipeline if it detects a hazard. Each stalls the corresponding stage and all stages preceding it for one more cycle. When a stage stalls, it keeps the instruction currently residing in it for another cycle and then restarts all stage activities exactly as in the non-stalled case.

The pipeline unit on which each instruction is to be executed is determined at run time by the instruction loader 18 when instructions are fetched from main memory.

The instruction loader 18 decodes prefetched instructions, tries to pack them into instruction pair entries and presents them to the dual-pipeline. If the instruction cache 22 is enabled (as discussed above), cacheable instructions can be stored in the instruction cache 22. In this case, an entry containing an instruction pair or a single instruction is also sent to the instruction cache 22 and stored there as a single cache entry. On instruction cache hits, stored instruction pairs are retrieved from the

instruction cache 22 and presented to the dual-pipeline for execution.

The instruction loader 18 attempts to pack instructions into pairs whenever possible. The packing of two instructions into one entry is possible only if the first instruction can be executed by integer pipeline unit 24 and both instructions are less than a preselected maximum length. If it is impossible to pack two instructions into a pair, then a single instruction is placed in Slot B.

Two instructions can be paired only when all of the following conditions hold: (1) both instructions are performance-critical core instructions, (2) the first instruction is executable by integer pipeline unit 24, and (3) the displacement and immediate fields in both instructions use short-encoding (short encoding for all instructions except the Branch instruction is 11 bits and 17 bits for the Conditional Branch and Branch and Link instructions).

Several instructions of the microprocessor 10 instruction set are restricted to run on integer pipeline unit 26 only. For example, because instruction pairs in the instruction cache 22 are tagged by the Slot A address, it is not useful to put a Branch instruction in Slot A since the corresponding Slot B instruction will not be accessible. Similarly, since there is a single arithmetic floating point pipe, it is not possible to execute two arithmetic floating point instructions in parallel. Restricting these instructions to integer pipeline unit 26 makes it possible to considerably simplify the dual-pipe data path design without hurting performance.

Integer unit 26 can execute any instructions in the microprocessor 10 instruction set.

The instruction loader 18 initiates instruction pairing upon an instruction cache miss, in which case it begins prefetching instructions into an instruction queue. In parallel, the instruction loader 18 examines the next instruction not yet removed from the instruction queue and attempts to pack it according to the following algorithm:

Step 1: Try to fit the next instruction into Slot A.

- (a) if the next instruction is not performance critical, then go to Step 5.
- (b) remove the next instruction from the instruction queue and tentatively place it in Slot A.
- (c) if the instruction is illegal for Slot A or if the instruction has an immediate/displacement field that cannot be represented in 11 bits, or if the instruction is not quad-word aligned, then go to Step 4.
- (d) otherwise, continue to Step 2.

Step 2: Try to fit the next instruction into Slot B.

- (a) if the next instruction is not performance-critical, or the next instruction has an encoded immediate/displacement field longer than 11 bits, or the next instruction is a branch with displacement longer than 17 bits, then go to Step 4.
- (b) otherwise, remove the next instruction from the

instruction queue, place it in Slot B and go to Step 3.

Step 3: Construct an instruction pair entry.

In this case, both Slot A and Slot B contain valid instructions and all pairing conditions are satisfied. Issue a pair entry and go to Step 1.

Step 4: Construct a single instruction entry.

In this case, Slot A contains an instruction which cannot be paired. Move this instruction to Slot B. If this instruction contains an immediate/displacement field longer than 17 bits, or it is a branch with displacement longer than 17 bits, and is not quad-word aligned, then replace it with UNDefined. Issue the entry and go to Step 1.

Step 5: Handle non-performance-critical instructions.

Remove the next instruction from the instruction queue and send it to the instruction emulator 20. When finished with this instruction, go to Step 1.

The just-described pairing algorithm packs two instructions whenever they can be held in a single instruction cache entry. However, these instructions may happen to be dependent, in which case they cannot be executed in parallel. The dependencies are detected by the execution processor 14.

## Claims

1. A processor that executes instructions, comprising:

a processor unit (14) with a plurality of functional units (24, 26, 28) for execution of instructions in parallel,

first retrieving means (18) for retrieving an encoded instruction from an external main memory;

decoding means (18,20) for decoding encoded instructions;

internal cache memory storage means (22) comprising a plurality of cache memory storage locations for storing instructions; and

second retrieving means for simultaneously retrieving a plurality of instructions from selected cache memory storage locations for execution by said functional units (24, 26, 28),

wherein

the internal cache memory storage means (22) is conceived for storing partially decoded instructions,

the encoded instructions retrieved from the main memory are partially decoded by said decoding means (18,20) before being stored in said internal cache memory storage means (22), and

said partially decoded instructions retrieved

from said internal cache memory storage means (22) are sent directly to the functional units (24, 26, 28) by said second retrieving means.

2. The processor according to claim 1, wherein each of said cache memory storage locations comprises a plurality of storage slots, each of said storage slots comprising means for storing a partially decoded instruction.

3. The processor according to claim 2, wherein said second retrieving means is capable of simultaneously retrieving a plurality of partially decoded instructions from said storage slots of a selected cache memory storage location for parallel execution by said plurality of functional units (24, 26, 28).

4. The processor according to one of claims 1 to 3, wherein each of said cache memory storage locations include means for storing auxiliary information indicative of whether the plurality of instructions stored in said slots of a cache memory storage location are independent such that the instructions may be executed in parallel or dependent such that the instructions must be executed sequentially.

5. The processor according to one of claims 1 to 4, wherein said external main memory is connected to the processor by means of a system bus, said system bus being connected to a bus interface unit (36) for retrieving encoded core instructions and encoded non-core instructions from said external main memory.

6. The processor according to one of claims 1 to 5, wherein the decoding means (18,20) comprises an instruction loader (18) for translating a first encoded core instruction to a first partially-decoded instruction and a second encoded core instruction to a second partially-decoded instruction further including means responsive to a received non-core instruction, wherein

said internal cache memory storage means (22) comprises said plurality of cache memory storage locations, each cache memory storage location comprising a plurality of storage slots, each of the storage slots comprising means for storing a decoded instruction; and means for simultaneously retrieving a plurality of decoded instructions from the storage slots of a selected cache memory storage location for parallel or sequential execution by the plurality of functional units, and wherein each of the cache memory storage locations includes means for storing auxiliary information indicative of whether the plurality of instructions

stored in the slots of a cache memory storage location are independent such that the instructions may be executed in parallel or dependent such that the instructions must be executed sequentially.

7. The processor according to one of claims 1 to 6, wherein the internal cache memory storage means (22) comprises a two-way set-associative organization.

8. A method of executing instructions in a processor, said processor comprising:

a processor unit (14) with a plurality of functional units (24, 26, 28) for execution of instructions in parallel;  
first retrieving means (18) for retrieving an encoded instruction from the external main memory;  
decoding means (18, 20) for decoding encoded instructions;  
internal cache memory storage means (22) comprising a plurality of cache memory storage locations for storing instructions; and  
second retrieving means for simultaneously retrieving a plurality of instructions from selected cache memory storage locations for execution by said functional units (24, 26, 28),

said method comprising the following steps:

- (a) retrieving encoded instructions from said external main memory;  
(b) partially decoding the instructions retrieved in retrieving step (a);  
(c) storing the instructions partially decoded in decoding step (b) in said internal cache memory storage means (22); and  
(d) retrieving the partially decoded instructions stored in storing step (c) for subsequent execution by said plurality of functional units,

wherein the steps (a) through (c) are performed before step (d) is performed, whereby partially decoded instructions are directly delivered to the functional units (24, 26, 28) thus enhancing the processing speed.

9. The method of claim 8, including the step of storing auxiliary information in the cache memory storage locations, the auxiliary information being indicative of whether the plurality of instructions stored in the slots of a cache memory storage location are independent such that the instructions may be executed in parallel, or dependent such that the instructions must be executed sequentially.

10. The method of claim 8 or claim 9, wherein said internal cache memory storage means comprises said plurality of cache memory storage locations, each cache memory storage location comprising a plurality of storage slots, each of the storage slots comprising means for storing a decoded instruction; and

simultaneously retrieving a plurality of decoded instructions from the storage slots of a selected cache memory storage location for parallel or sequential execution by the plurality of functional units, and including the step of storing auxiliary information in the cache memory storage locations, the auxiliary information being indicative of whether the plurality of instructions stored in the slots of a cache memory storage location are independent such that the instructions may be executed in parallel, or dependent such that the instructions must be executed sequentially.

#### Patentansprüche

1. Prozessor, der Befehle ausführt, mit:

einer Prozessoreinheit (14), die mehrere funktionale Einheiten (24, 26, 28) zum parallelen Ausführen von Befehlen enthält,  
einer ersten Wiedergewinnungseinrichtung (18) zum Wiedergewinnen eines codierten Befehls von einem externen Hauptspeicher;  
einer Decodierungseinrichtung (18, 20) zum Decodieren codierter Befehle;  
einer internen Cache-Speichereinrichtung (22), die mehrere Cache-Speicherplätze zum Speichern von Befehlen enthält; und  
einer zweiten Wiedergewinnungseinrichtung zum gleichzeitigen Wiedergewinnen mehrerer Befehle von ausgewählten Cache-Speicherplätzen für deren Ausführung durch die funktionalen Einheiten (24, 26, 28),

wobei

die interne Cache-Speichereinrichtung (22) so beschaffen ist, daß sie teilweise decodierte Befehle speichert,  
die codierten Befehle, die aus dem Hauptspeicher wiedergewonnen werden, durch die Decodierungseinrichtung (18, 20) teilweise decodiert werden, bevor sie in der internen Cache-Speichereinrichtung (22) gespeichert werden, und  
die teilweise decodierten Befehle, die von der internen Cache-Speichereinrichtung (22) wiedergewonnen werden, von der zweiten Wiedergewinnungseinrichtung direkt zu den funktionalen Einheiten (24, 26, 28) geschickt werden.

2. Prozessor nach Anspruch 1, wobei jeder der Cache-Speicherplätze mehrere Speicherschlitzte enthält, wovon jeder eine Einrichtung zum Speichern eines teilweise decodierten Befehls enthält. 5
3. Prozessor nach Anspruch 2, wobei die zweite Wiedergewinnungseinrichtung gleichzeitig mehrere teilweise decodierte Befehle von den Speicherschlitzten eines ausgewählten Cache-Speicherplatzes wiedergewinnen kann, damit sie von den mehreren funktionalen Einheiten (24, 26, 28) parallel ausgeführt werden können. 10
4. Prozessor nach einem der Ansprüche 1 bis 3, wobei jeder der Cache-Speicherplätze eine Einrichtung zum Speichern von Hilfsinformationen enthält, die angeben, ob die mehreren in den Schlitzten eines Cache-Speicherplatzes gespeicherten Befehle unabhängig sind, so daß die Befehle parallel ausgeführt werden können, oder abhängig sind, so daß die Befehle sequentiell ausgeführt werden müssen. 15 20
5. Prozessor nach einem der Ansprüche 1 bis 4, wobei der externe Hauptspeicher an den Prozessor durch einen Systembus angeschlossen ist, der an eine Busschnittstelleneinheit (36) zum Wiedergewinnen codierter Kernbefehle und codierter Nichtkern-Befehle vom externen Hauptspeicher angeschlossen ist. 25 30
6. Prozessor nach einem der Ansprüche 1 bis 5, wobei die Decodierungseinrichtung (18, 20) eine Befehls-ladeeinrichtung (18) zum Übersetzen eines ersten codierten Kernbefehls in einen ersten teilweise decodierten Befehl und eines zweiten codierten Kernbefehls in einen zweiten teilweise decodierten Befehl enthält und ferner eine Einrichtung enthält, die auf einen empfangenen Nichtkern-Befehl anspricht, wobei 35 40
 

die interne Cache-Speichereinrichtung (22) die mehreren Cache-Speicherplätze umfaßt, wobei jeder Cache-Speicherplatz mehrere Speicherschlitzte enthält, wovon jeder eine Einrichtung zum Speichern eines decodierten Befehls enthält; und 45

eine Einrichtung zum gleichzeitigen Wiedergewinnen mehrerer decodierter Befehle aus den Speicherschlitzten eines ausgewählten Cache-Speicherplatzes zum parallelen oder sequentiellen Ausführen durch die mehreren funktionalen Einheiten, und wobei 50

jeder der Cache-Speicherplätze eine Einrichtung zum Speichern von Hilfsinformationen enthält, die angeben, ob die mehreren Befehle, die in den Schlitzten eines Cache-Speicherplatzes gespeichert sind, unabhängig sind, so daß die Befehle parallel ausgeführt werden können, 55

oder abhängig sind, so daß die Befehle sequentiell ausgeführt werden müssen.

7. Prozessor nach einem der Ansprüche 1 bis 6, wobei die interne Cache-Speichereinrichtung (22) eine mengenassoziative Zweiwege-Organisation enthält.
8. Verfahren zum Ausführen von Befehlen in einem Prozessor, wobei der Prozessor enthält:
 

eine Prozessoreinheit (14), die mehrere funktionale Einheiten (24, 26, 28) zum parallelen Ausführen von Befehlen enthält;

eine erste Wiedergewinnungseinrichtung (18) zum Wiedergewinnen eines codierten Befehls vom externen Hauptspeicher;

eine Decodierungseinrichtung (18, 20) zum Decodieren codierter Befehle;

eine interne Cache-Speichereinrichtung (22), die mehrere Cache-Speicherplätze zum Speichern von Befehlen enthält; und

eine zweite Wiedergewinnungseinrichtung zum gleichzeitigen Wiedergewinnen mehrerer Befehle von ausgewählten Cache-Speicherplätzen für die Ausführung durch die funktionalen Einheiten (24, 26, 28),

wobei das Verfahren die folgenden Schritte enthält:

  - (a) Wiedergewinnen codierter Befehle vom externen Hauptspeicher;
  - (b) teilweises Decodieren der im Wiedergewinnungsschritt (a) wiedergewonnen Befehle;
  - (c) Speichern der im Decodierungsschritt (b) teilweise decodierten Befehle in der internen Cache-Speichereinrichtung (22); und
  - (d) Wiedergewinnen der im Speicherschritt (c) gespeicherten teilweise decodierten Befehle für eine nachfolgende Ausführung durch die mehreren funktionalen Einheiten,

wobei die Schritte (a) bis (c) vor der Ausführung des Schrittes (d) ausgeführt werden, wobei teilweise decodierte Befehle direkt zu den funktionalen Einheiten (24, 26, 28) geschickt werden, wodurch die Verarbeitungsgeschwindigkeit erhöht wird.
9. Verfahren nach Anspruch 8, mit dem Schritt des Speicherns von Hilfsinformationen in den Cache-Speicherplätzen, wobei die Hilfsinformationen angeben, ob die mehreren in den Schlitzten eines Cache-Speicherplatzes gespeicherten Befehle unabhängig sind, so daß die Befehle parallel ausgeführt werden können, oder abhängig sind, so daß die Befehle sequentiell ausgeführt werden müssen.
10. Verfahren nach Anspruch 8 oder 9, bei dem die in-

terne Cache-Speichereinrichtung die mehreren Cache-Speicherplätze enthält, wobei jeder Cache-Speicherplatz mehrere Speicherschlitzte enthält, wovon jeder eine Einrichtung zum Speichern eines decodierten Befehls enthält; und

mit dem Schritt des gleichzeitigen Wiedergewinnens mehrerer decodierter Befehle aus den Speicherschlitzten eines ausgewählten Cache-Speicherplatzes für eine parallele oder sequentielle Ausführung durch die mehreren funktionalen Einheiten, und mit dem Schritt des Speicherns von Hilfsinformationen in den Cache-Speicherplätzen, wobei die Hilfsinformationen angeben, ob die mehreren in den Schlitzten eines Cache-Speicherplatzes gespeicherten Befehle unabhängig sind, so daß die Befehle parallel ausgeführt werden können, oder abhängig sind, so daß die Befehle sequentiell ausgeführt werden müssen.

#### Revendications

1. Processeur qui exécute des instructions, comportant :

une unité de processeur (14) pourvue d'une pluralité d'unités fonctionnelles (24, 26, 28) pour l'exécution d'instructions en parallèle, des premiers moyens d'extraction (18) pour extraire une instruction codée d'une mémoire principale externe; des moyens de décodage (18, 20) pour décoder des instructions codées; des moyens de mémorisation à antémémoire interne (22) comportant une pluralité d'emplacements de mémorisation en antémémoire pour mémoriser des instructions; et des seconds moyens d'extraction pour extraire simultanément une pluralité d'instructions d'emplacements sélectionnés de mémorisation en antémémoire pour leur exécution par lesdites unités fonctionnelles (24, 26, 28), dans lequel les moyens de mémorisation à antémémoire interne (22) sont conçus pour mémoriser des instructions partiellement décodées, les instructions codées extraites de la mémoire principale sont partiellement décodées par lesdits moyens de décodage (18, 20) avant d'être mémorisées dans lesdits moyens de mémorisation à antémémoire interne (22), et lesdites instructions partiellement décodées extraites desdits moyens de mémorisation à antémémoire interne (22) sont envoyées directement aux unités fonctionnelles (24, 26, 28) par lesdits seconds moyens d'extraction.

2. Processeur selon la revendication 1, dans lequel chacun desdits emplacements de mémorisation en antémémoire comprend une pluralité de créneaux de mémorisation, chacun desdits créneaux de mémorisation comportant des moyens pour mémoriser une instruction partiellement décodée.

3. Processeur selon la revendication 2, dans lequel lesdits seconds moyens d'extraction sont capables d'extraire simultanément une pluralité d'instructions partiellement décodées desdits créneaux de mémorisation d'un emplacement sélectionné de mémorisation en antémémoire pour leur exécution en parallèle par ladite pluralité d'unités fonctionnelles (24, 26, 28).

4. Processeur selon l'une des revendications 1 à 3, dans lequel chacun desdits emplacements de mémorisation en antémémoire comporte des moyens pour mémoriser des informations auxiliaires indiquant si les plusieurs instructions mémorisées dans lesdits créneaux d'un emplacement de mémorisation en antémémoire sont indépendantes, de sorte que les instructions peuvent être exécutées en parallèle, ou dépendantes de sorte que les instructions doivent être exécutées en séquence.

5. Processeur selon l'une des revendications 1 à 4, dans lequel ladite mémoire principale externe est connectée au processeur au moyen d'un bus de système, ledit bus de système étant connecté à une unité d'interface de bus (36) pour extraire des instructions codées critiques et des instructions codées non critiques de ladite mémoire principale externe.

6. Processeur selon l'une des revendications 1 à 5, dans lequel les moyens de décodage (18, 20) comportent un chargeur d'instructions (18) pour traduire une première instruction codée critique en une première instruction partiellement décodée et une seconde instruction codée critique en une seconde instruction partiellement décodée, comportant en outre des moyens répondant à la réception d'une instruction non critique, dans lequel

lesdits moyens de mémorisation à antémémoire interne (22) comportent ladite pluralité d'emplacements de mémorisation en antémémoire, chaque emplacement de mémorisation en antémémoire comprenant une pluralité de créneaux de mémorisation, chacun des créneaux de mémorisation comportant des moyens pour mémoriser une instruction décodée; et des moyens pour extraire simultanément une pluralité d'instructions décodées des créneaux de mémorisation d'un emplacement sélectionné de mémorisation en antémémoire pour leur

exécution en parallèle ou en séquence par la pluralité d'unités fonctionnelles, et dans lequel chacun des emplacements de mémorisation en antémémoire comporte des moyens pour mémoriser des informations auxiliaires indiquant si les plusieurs instructions mémorisées dans les créneaux d'un emplacement de mémorisation en antémémoire sont indépendantes, de sorte que les instructions peuvent être exécutées en parallèle, ou dépendantes de sorte que les instructions doivent être exécutées en séquence.

7. Processeur selon l'une des revendications 1 à 6, dans lequel les moyens de mémorisation à antémémoire interne (22) comportent une organisation associative bidirectionnelle.

8. Procédé d'exécution d'instructions dans un processeur, ledit processeur comportant :

une unité de processeur (14) pourvue d'une pluralité d'unités fonctionnelles (24, 26, 28) pour l'exécution d'instructions en parallèle;  
des premiers moyens d'extraction (18) pour extraire une instruction codée de la mémoire principale externe;  
des moyens de décodage (18, 20) pour décoder des instructions codées ;  
des moyens de mémorisation à antémémoire interne (22) comportant une pluralité d'emplacements de mémorisation en antémémoire pour mémoriser des instructions; et  
des seconds moyens d'extraction pour extraire simultanément une pluralité d'instructions d'emplacements sélectionnés de mémorisation en antémémoire pour leur exécution par lesdites unités fonctionnelles (24, 26, 28),

ledit procédé comportant les étapes suivantes :

(a) extraire des instructions codées de ladite mémoire principale externe;  
(b) décoder partiellement les instructions extraites dans l'étape d'extraction (a);  
(c) mémoriser les instructions partiellement décodées dans l'étape de décodage (b) dans lesdits moyens de mémorisation à antémémoire interne (22); et  
(d) extraire les instructions partiellement décodées mémorisées dans l'étape de mémorisation (c) pour leur exécution subséquente par ladite pluralité d'unités fonctionnelles,

dans lequel les étapes (a) à (c) sont effectuées avant que l'étape (d) soit effectuée, de sorte que des instructions partiellement décodées sont délivrées directement aux unités fonctionnelles (24,

26, 28), améliorant ainsi la vitesse de traitement.

9. Procédé selon la revendication 8, comportant l'étape de mémoriser des informations auxiliaires dans les emplacements de mémorisation en antémémoire, les informations auxiliaires indiquant si les plusieurs instructions mémorisées dans les créneaux d'un emplacement de mémorisation en antémémoire sont indépendantes, de sorte que les instructions peuvent être exécutées en parallèle, ou dépendantes de sorte que les instructions doivent être exécutées en séquence.

10. Procédé selon la revendication 8 ou 9, dans lequel lesdits moyens de mémorisation à antémémoire interne (22) comportent ladite pluralité d'emplacements de mémorisation en antémémoire, chaque emplacement de mémorisation en antémémoire comprenant une pluralité de créneaux de mémorisation, chacun des créneaux de mémorisation comportant des moyens pour mémoriser une instruction décodée; et

on extrait simultanément une pluralité d'instructions décodées des créneaux de mémorisation d'un emplacement sélectionné de mémorisation en antémémoire pour leur exécution en parallèle ou en séquence par la pluralité d'unités fonctionnelles, et comportant l'étape de mémorisation d'informations auxiliaires dans les emplacements de mémorisation en antémémoire, les informations auxiliaires indiquant si les plusieurs instructions mémorisées dans les créneaux d'un emplacement de mémorisation en antémémoire sont indépendantes, de sorte que les instructions peuvent être exécutées en parallèle, ou dépendantes de sorte que les instructions doivent être exécutées en séquence.

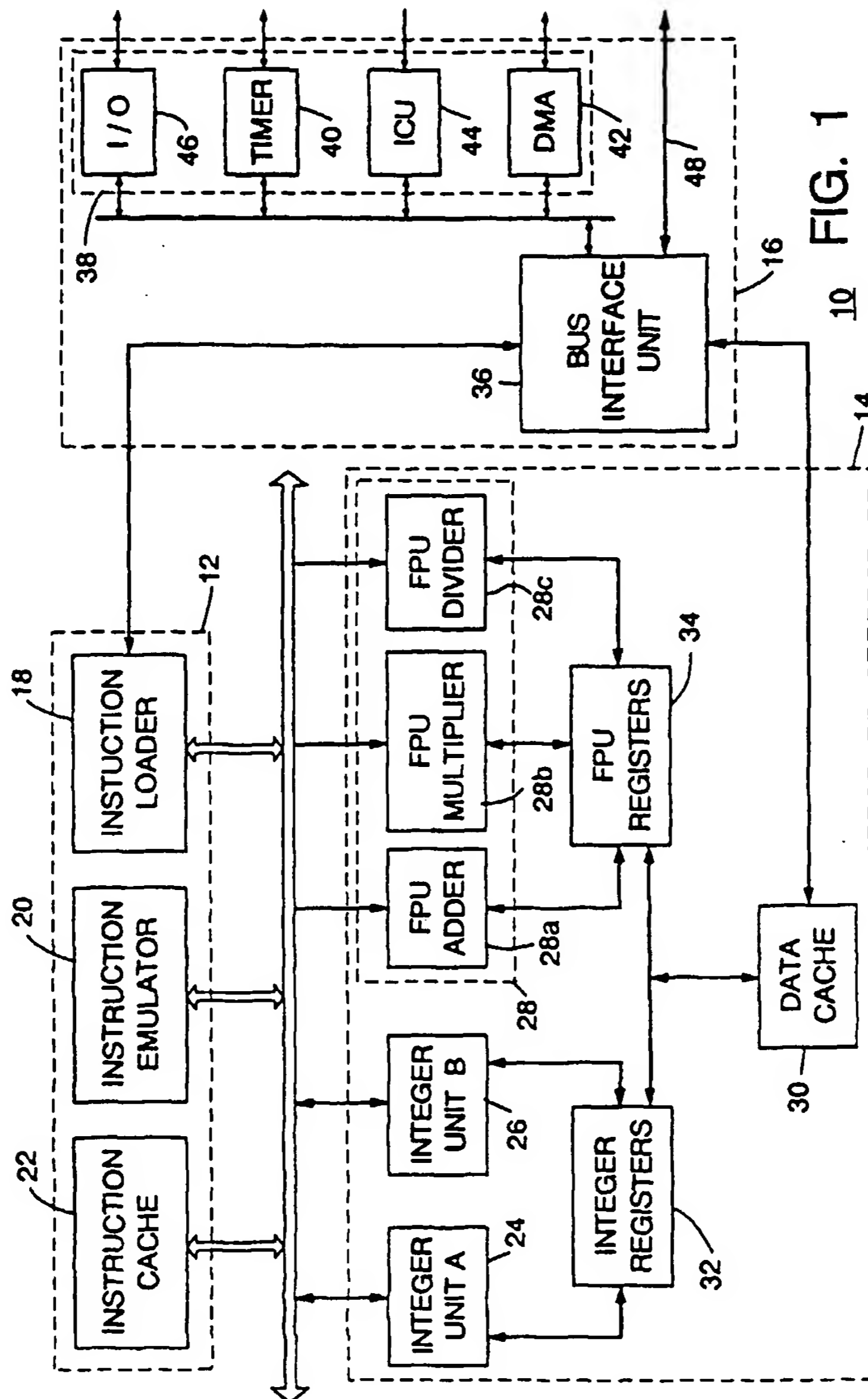


FIG. 1

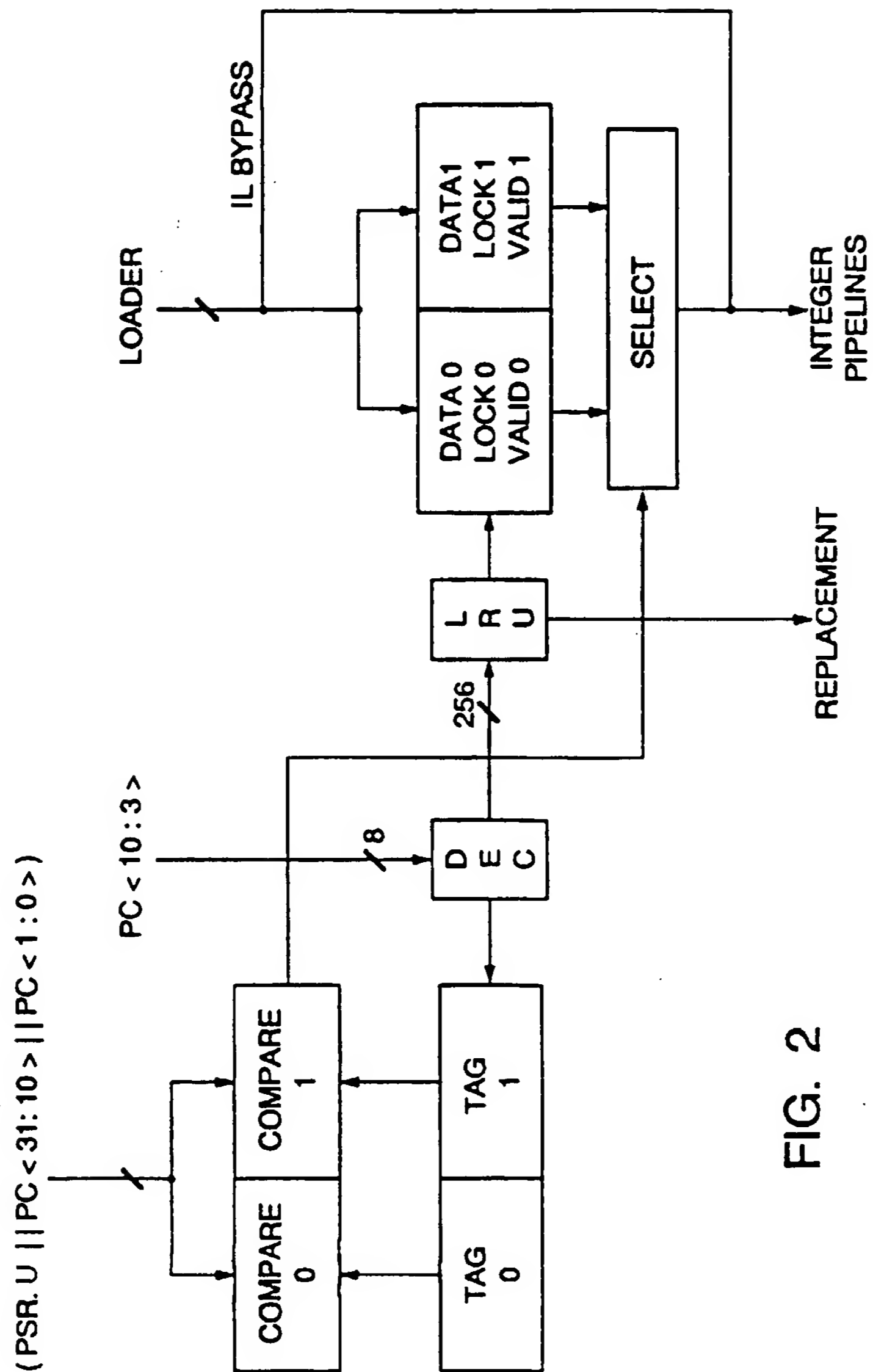
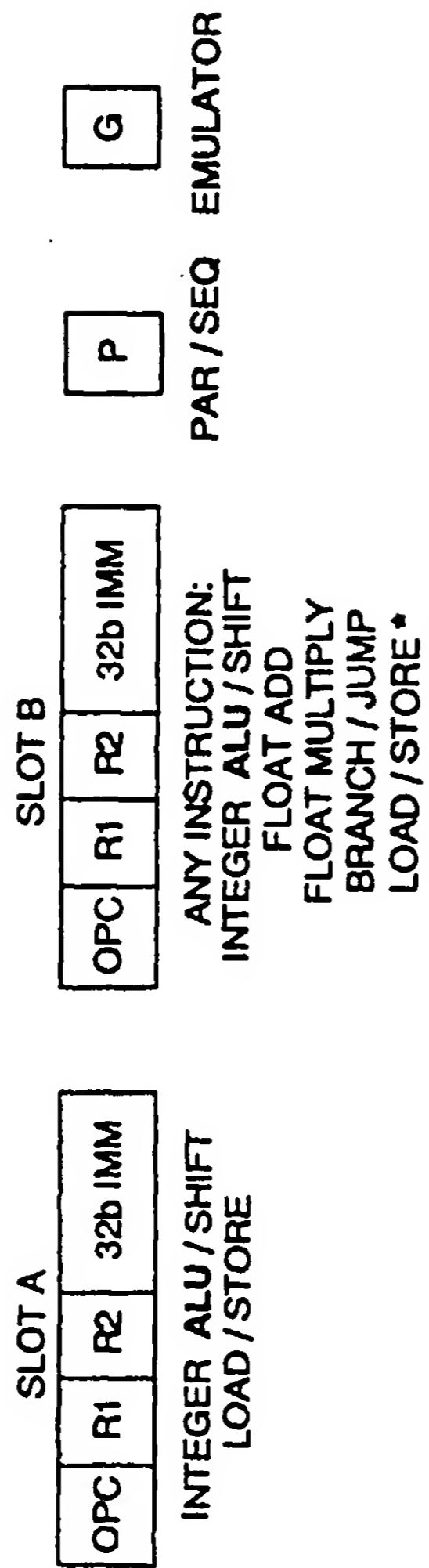


FIG. 2



\* LOAD / STORE IN SLOT B ONLY FOR SEQUENTIAL EXECUTION

FIG. 3

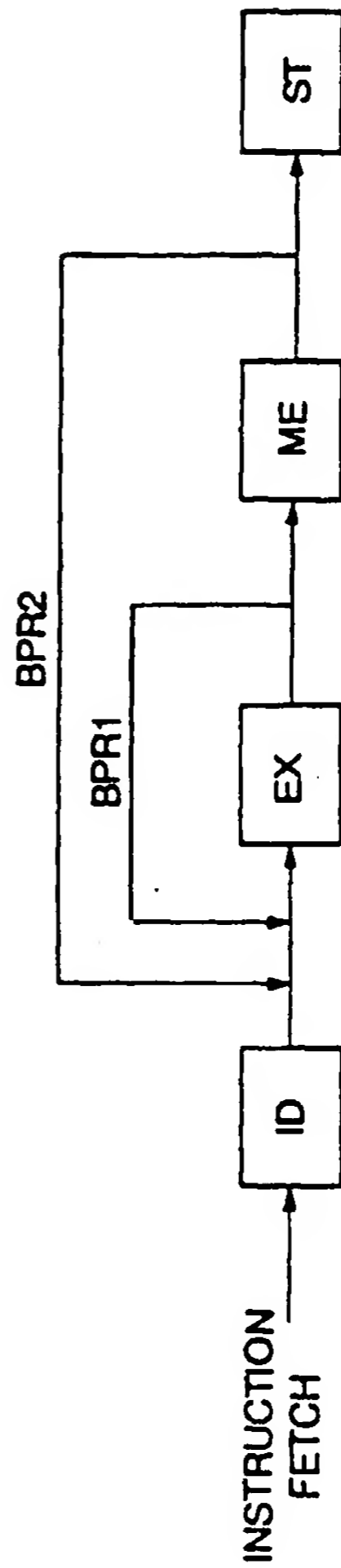


FIG. 4